

Basic Models

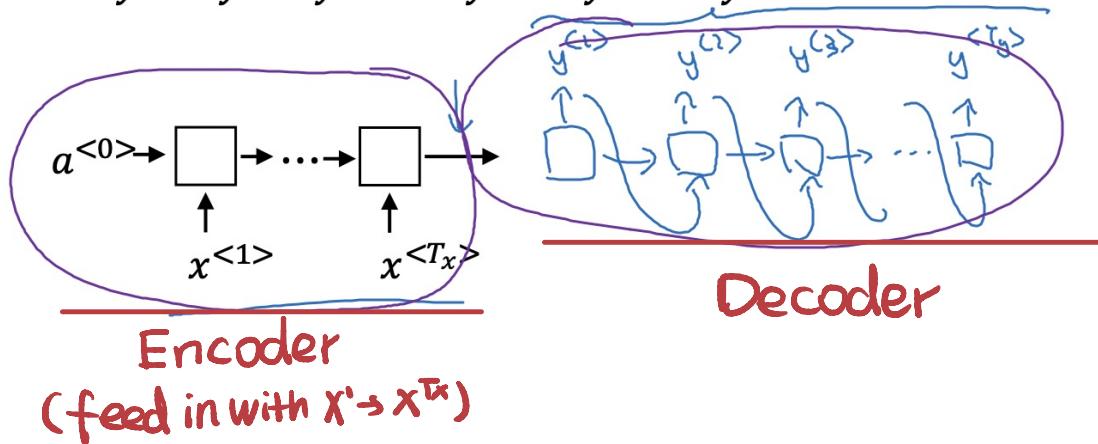
Sequence-to-sequence model :

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

Jane visite l'Afrique en septembre

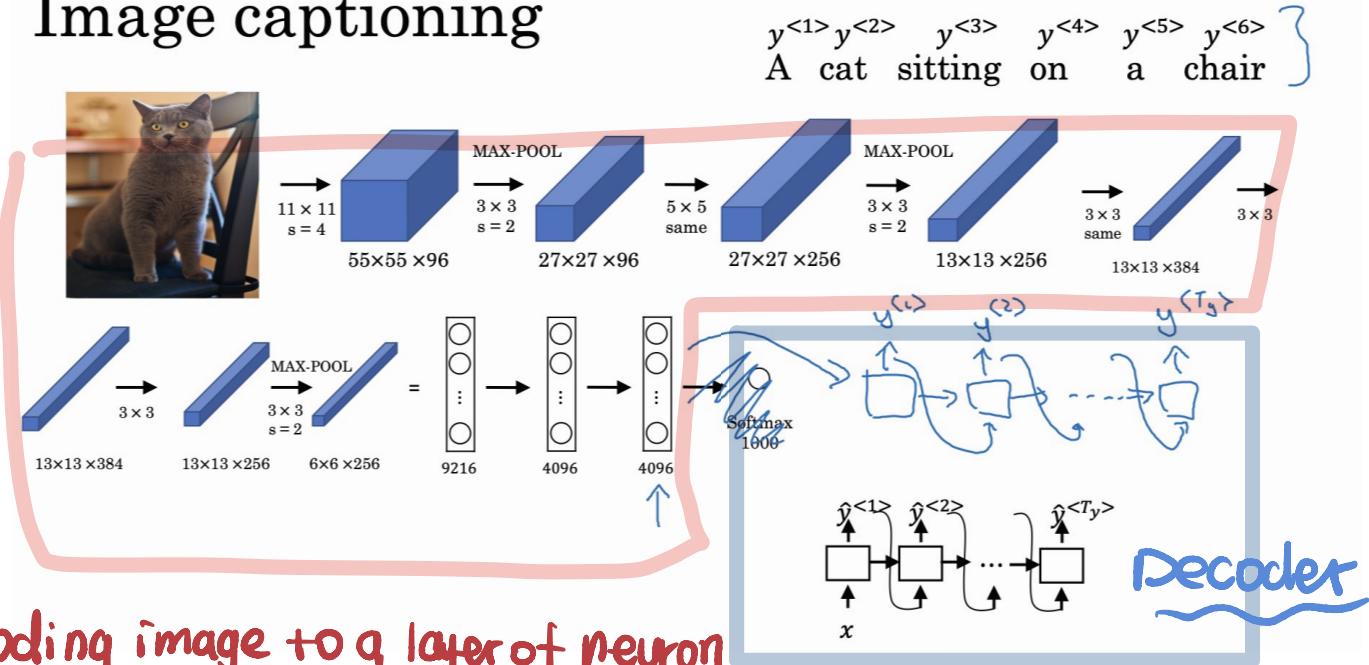
→ Jane is visiting Africa in September.

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$

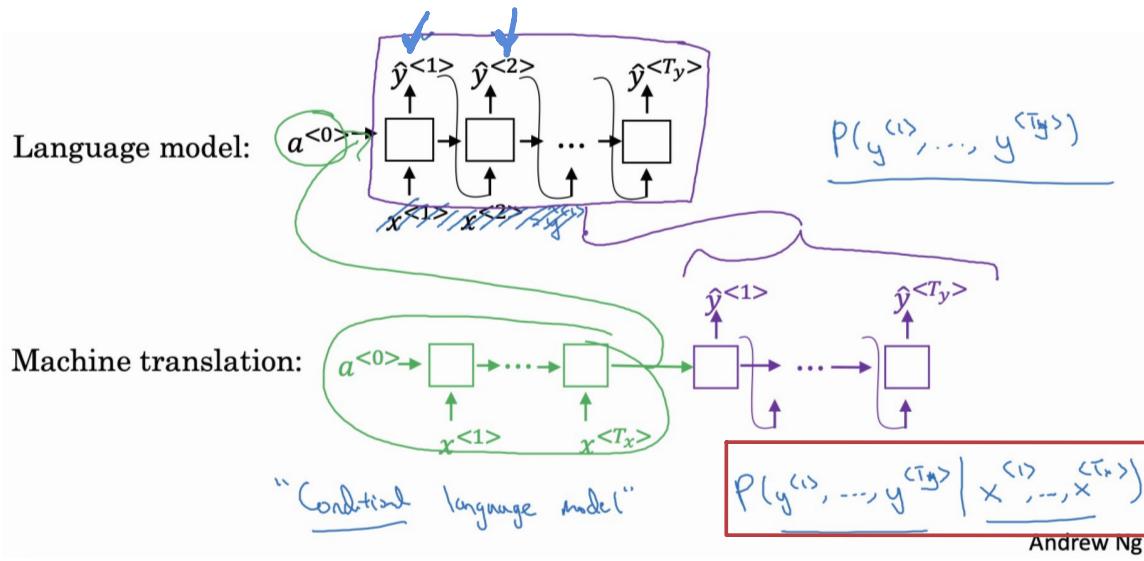


Same encoder-decoder technique can be used to create Image Captioning Model.

Image captioning



Picking the Most Likely Sentence



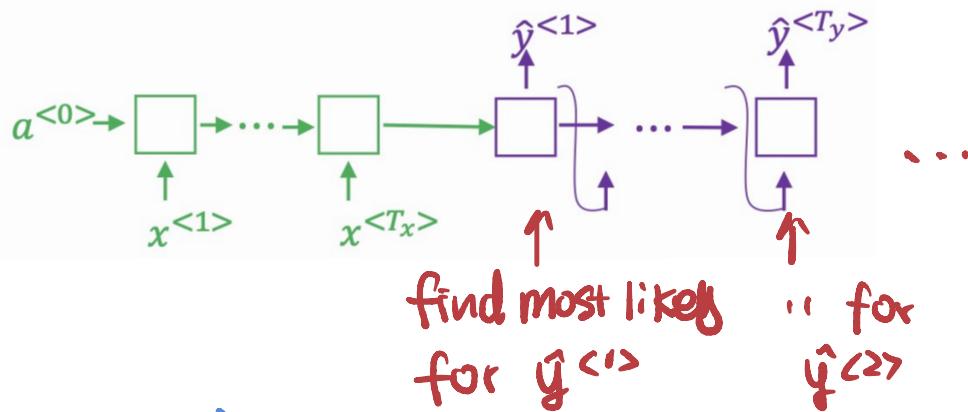
Another difference :

conditioning Probability
given x

Andrew Ng

With language model, we want to generate output randomly from the predicted distribution
In machine translation: we generate best translation (most likely translation)

Naive Approach: greedy search



Problem: → Jane is visiting Africa in September
 → Jane is going to be visiting Africa in September

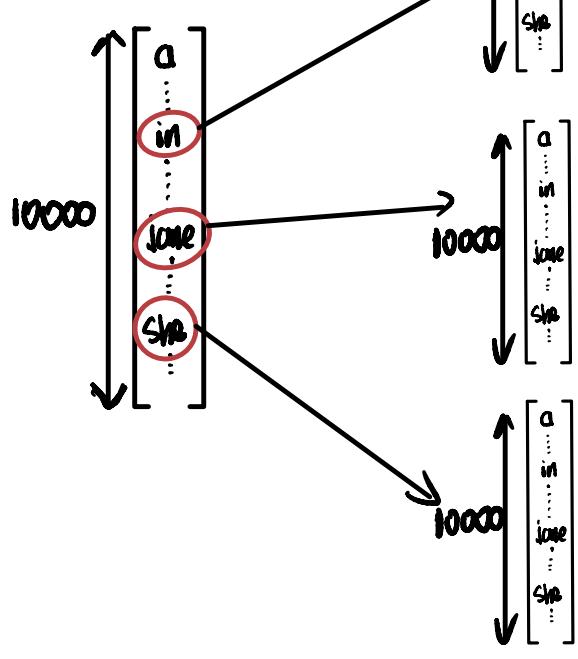
Even though sentence 1) is a better translation. Since "going" is used more frequently used in human text, sentence 2) is more likely to be output by greedy

→ Better Solution

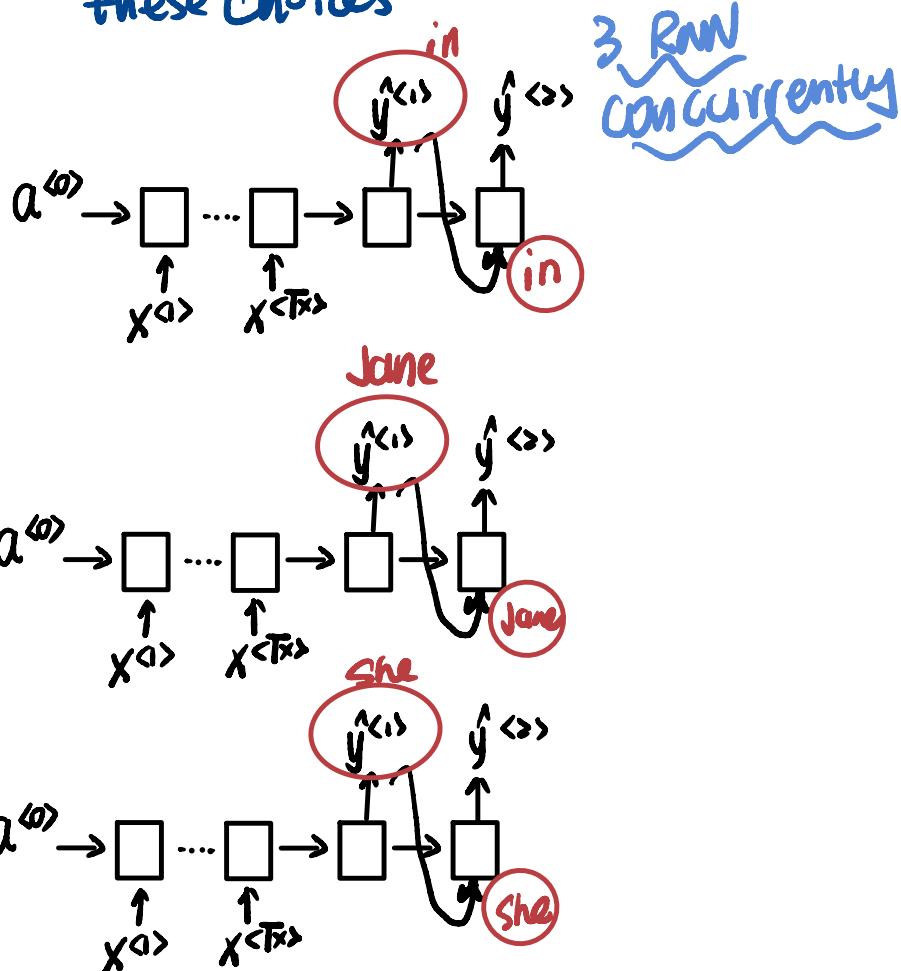
Beam Search

B=3 (Beam Width)

Step 1: choose top 3 words as 1st word choice



Step 2, Combine each choice from Step 1 with whole Vocabulary to form $3 \times 10,000 = 30,000$ choices, choose top 3 - Beam width from these choices



$$\rightarrow P(y^{<1>} | x, y^{<2>})$$

$$= P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

length Normalization

$$\operatorname{argmax}_y \prod_{t=1}^T P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>})$$

Not very numerically stable, prone to under-flow, since we are multiplying many <1 numbers.

instead, we use: $\operatorname{argmax}_y \sum_{t=1}^T \log P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>})$

$$\frac{1}{T_{ya}} \sum_{y=1}^T \log P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>}) \quad \alpha = 0.7, 1, \dots$$

hyperparameter

weighting out the undesirable negative effect
of low possibility assigned to relatively long sentences.

Beam Width B ?

large B : better result, slower

small B : worse result, faster.

$B=10$

Common

100,

Production

1000 → 3000

theoretical papers

really big.

- Unlike exact search algorithms, such as BFS and DFS, Beam Search runs faster but is not guaranteed to find exact maximum for $\arg \max_y P(y|x)$

Error Analysis on Beam Search

Error analysis on beam search

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y})

Case 1: $P(y^*|x) > P(\hat{y}|x)$ ←

$\arg \max_y P(y|x)$

Beam search chose \hat{y} . But y^* attains higher $P(y|x)$.

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \leq P(\hat{y}|x)$ ←

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) < P(\hat{y}|x)$.

Conclusion: RNN model is at fault.

$P(y^*|x)$
 $P(\hat{y}|x)$

computed thru forward propagation through RNN

Error Analysis process: Pick multiple false preds. Andrew Ng
Perform the above analysis → What fraction of error → beam RNN model.

Bleu Score

$$\text{Bleu Score on unigram} = \frac{\sum_{\text{Unigrame } j} \text{Count}_{\text{clip}}(\text{Unigram})}{\sum_{\text{Unigrame } j} \text{count}(\text{Unigram})}$$

$$P_2 = \frac{\sum_{n\text{-gram } j} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram } j} \text{count}(n\text{-gram})}$$

P_n = Bleu score on n -grams only P_1, P_2, P_3, P_4

combined Bleu score : $\tilde{BP} \exp \left(\frac{1}{4} \sum_{n=1}^4 P_n \right)$
brevity penalty

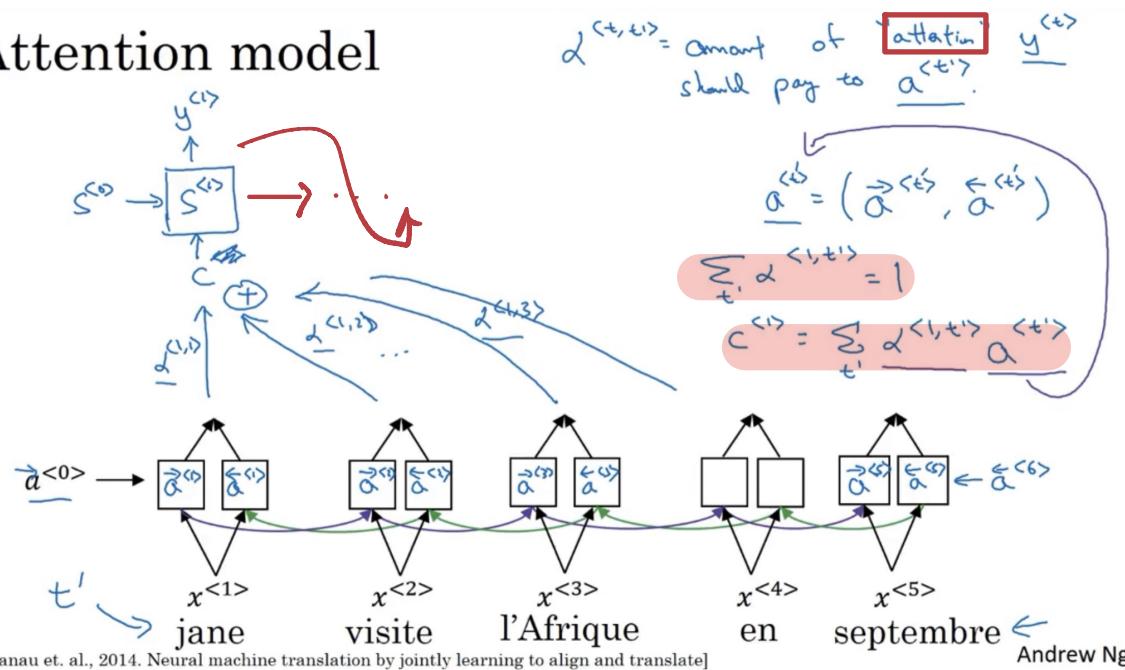
$$\tilde{BP} = \begin{cases} 1 & \text{if MT-output-length} > \text{reference-length} \\ \exp(1 - \text{ref-length} / \text{MT-output-length}) & \end{cases}$$

Attention Model

Intuition: When human translates, we won't read whole sentence and output entire translation in one iteration. Rather, we read in and output translation part by part at the same time.

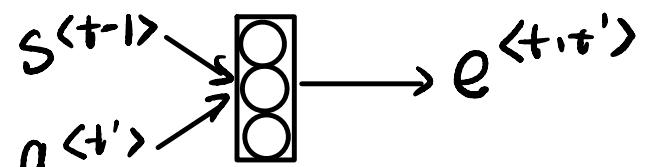
Similarly, neural networks are also bad at processing large sequence of input at once. Thus we designed attention model which can figure out how much attention we should pay to nearby words when translating current word.

Attention model



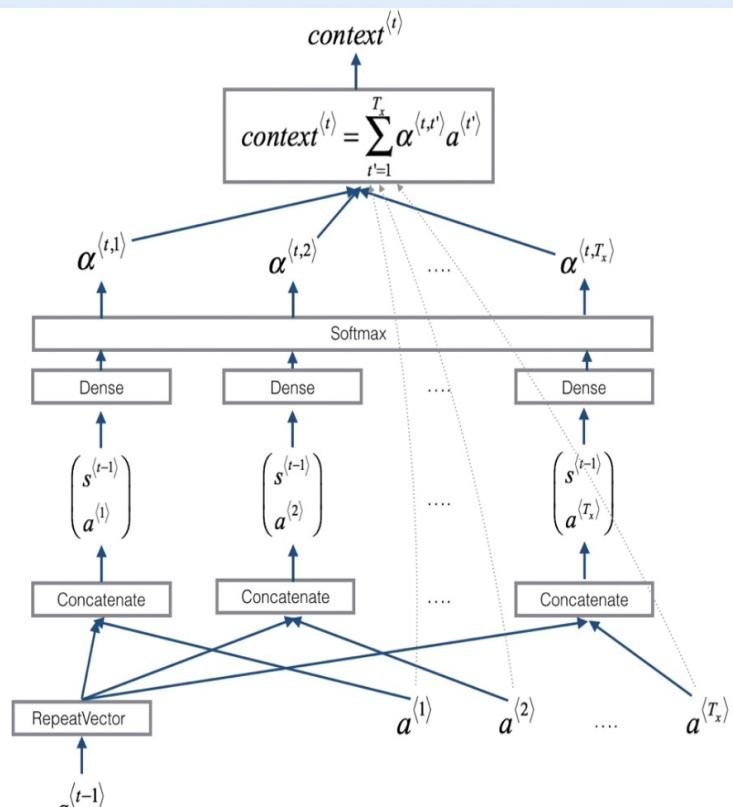
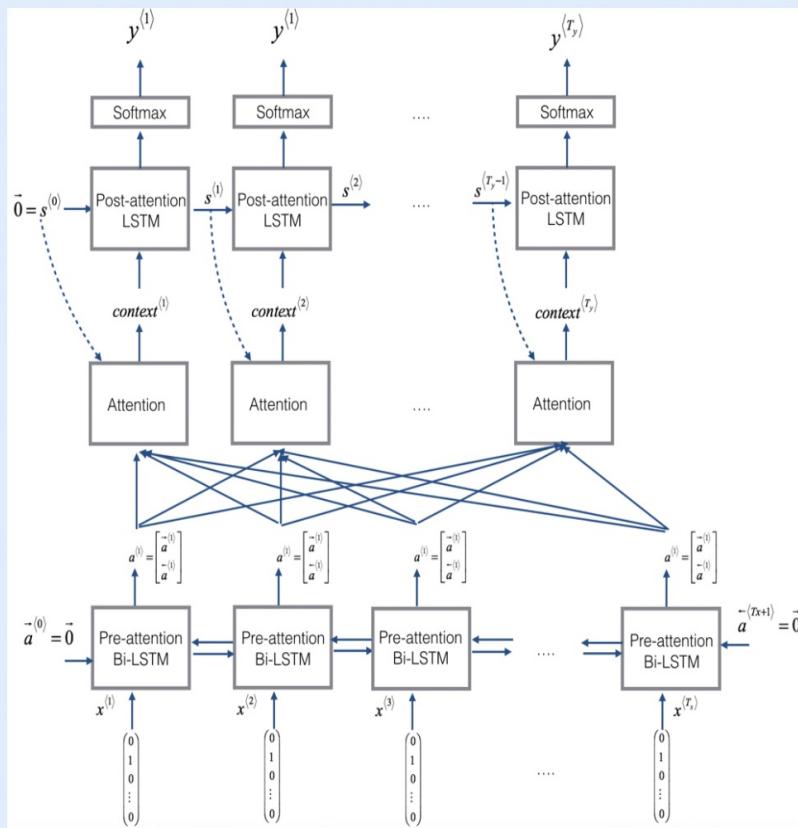
$\alpha^{(t,t')}$ = amount of attention $y^{(t)}$ should pay to $a^{(t')}$

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^T \exp(e^{(t,t')})}$$



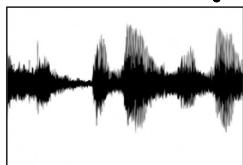
Attention depends on last translation and coming hidden activation from encoder.

- The diagram on the left shows the attention model.
- The diagram on the right shows what one "attention" step does to calculate the attention variables $\alpha^{(t,t')}$.
- The attention variables $\alpha^{(t,t')}$ are used to compute the context variable $context^{(t)}$ for each timestep in the output ($t = 1, \dots, T_y$).

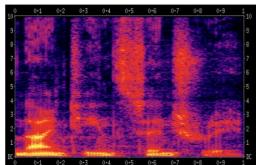


Speech Recognition

x
audio clip \longrightarrow y
transcript

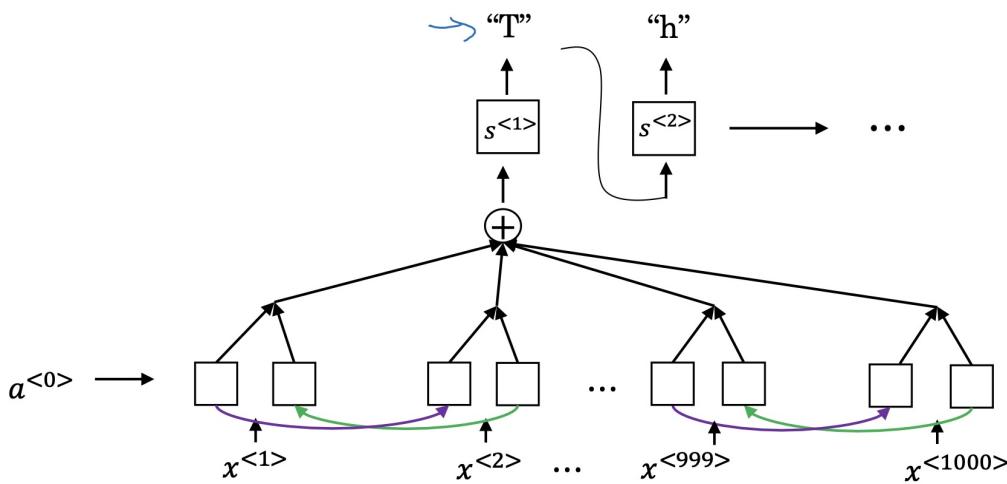


↓ possible preprocessing



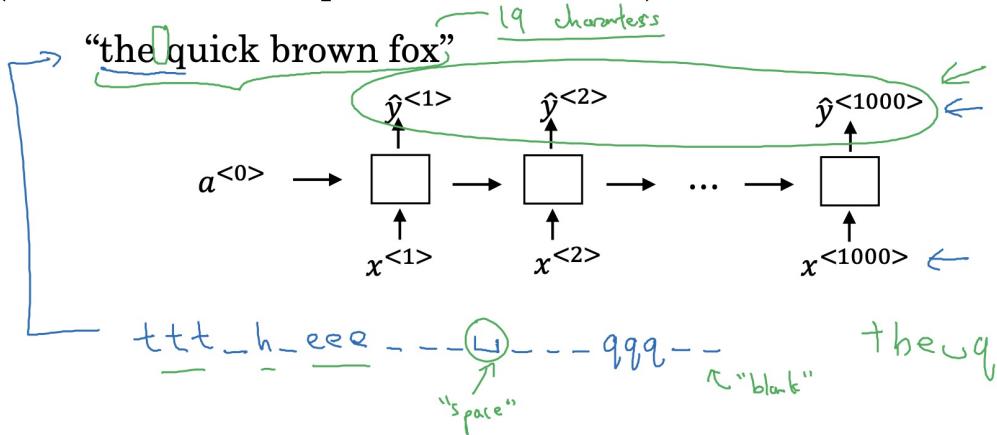
(spectrographs)

Attention model for speech recognition



CTC cost for speech recognition

(Connectionist temporal classification)



Basic rule: collapse repeated characters not separated by "blank"

Trigger word Detection

